

文本即数据实验报告

李浩然 2025 年 11 月 29 日

1 实验目的

- 深入理解 LC-3 汇编语言的指令格式、寻址方式和汇编器的工作原理。
- 掌握模拟汇编器的实现方法，包括符号表构建、指令翻译和地址计算。
- 能够使用 Python 语言实现一个简易的 LC-3 汇编器，完成从汇编指令到机器码的转换。
- 验证汇编器对各类 LC-3 指令（包括伪指令、陷阱指令、分支指令等）的处理能力。

2 实验原理

2.1 LC-3 汇编语言基础

LC-3 (Little Computer 3) 是一种简化的教学用计算机体系结构，其汇编语言包含**操作码 (opcode)**、**寄存器操作数**、**立即数**和**标签**等元素。指令长度固定为 16 位，主要分为运算类、访存类、控制转移类和陷阱指令类，同时支持.`ORIG`、`END`、`FILL`、`BLKW`、`STRINGZ` 等伪指令用于定义程序起始地址、数据空间和字符串常量。

2.2 汇编器工作原理

汇编器的核心是将汇编指令转换为对应的机器码，通常采用**两趟扫描法**实现：

- 第一趟扫描**: 遍历汇编源代码，构建**符号表 (标签-地址映射)**，同时计算程序计数器 (PC) 的值，处理伪指令对地址的占用，记录每条指令的地址和原始指令信息。
- 第二趟扫描**: 根据符号表和指令格式，将每条汇编指令翻译为 16 位机器码，处理地址偏移量计算、寄存器编码和立即数转换等细节。

3 实验环境

- 操作系统: Alpine Linux 3.22.2
- 编程语言: Python 3.14.0 (main, Oct 8 2025, 23:12:31) [GCC 14.2.0]
- 开发工具: VS Code, Ubuntu
- 测试工具: 文本编辑器（用于编写测试用例）、终端（用于运行汇编器并验证输出）

4 实验内容与实现

4.1 汇编器整体架构

本次我实现的 LC-3 汇编器通过 Python 语言编写，主要分为输入处理、第一趟扫描（符号表构建）、第二趟扫描（指令翻译）和输出生成四个模块，核心代码见 `main.py`。

4.2 核心函数实现

4.2.1 数据类型转换函数

- `to_int(x)`: 将十六进制（以 `x` 开头）、十进制（以 `#` 开头）或普通数字字符串转换为整数。
- `to_bin16(x)`: 将整数转换为 16 位二进制字符串，确保结果为 16 位补码形式。
- `reg(r)`: 从寄存器字符串（如 `R0`）中提取寄存器编号（0-7）。

4.2.2 地址偏移计算函数

`pcoff(pc_from, target, bits)`: 计算基于 PC 的偏移量，支持 9 位、11 位等不同长度的偏移量计算，处理负数偏移的补码转换。

4.2.3 指令处理函数

`handle_instruction(op, parts_raw, parts_up, pc)`: 根据操作码类型，分别处理运算指令（ADD/AND/NOT）、访存指令（LD/LDI/LDR/ST/STI/STR）、控制转移指令（BR/JMP/JSR/JSRR）、陷阱指令（TRAP）和伪指令（.FILL/.BLKW/.STRINGZ），生成对应的 16 位机器码。

4.3 两趟扫描实现

- 第一趟扫描：
 - 处理`.ORIG` 伪指令，初始化程序计数器（PC）。
 - 识别标签并构建符号表，记录标签对应的内存地址。
 - 处理`.BLKW`（分配连续内存空间）和`.STRINGZ`（存储字符串，末尾加 0）伪指令，更新 PC 值。
 - 记录每条指令的地址、原始指令和大写指令信息，供第二趟扫描使用。
- 第二胎扫描：
 - 遍历第一胎扫描记录的指令信息，根据操作码调用 `handle_instruction` 函数生成机器码。
 - 将生成的机器码按顺序存储，最终输出所有 16 位二进制机器码。

4.4 核心代码展示

4.4.1 main.py 核心代码

```

1 import sys
2 import re
3
4 def to_int(x):
5     """Convert string to integer based on prefix (hex or decimal)."""
6     x = x.lower()
7     if x.startswith("x"):
8         return int(x[1:], 16)
9     if x.startswith("#"):
10        return int(x[1:], 10)
11    return int(x)
12
13 def to_bin16(x):
14     """Convert integer to 16-bit binary string."""
15     return format(x & 0xFFFF, "016b")
16
17 def reg(r):
18     """Extract register number from the register string (e.g., R0 -> 0)."""
19     return int(r.strip(",")[-1])
20
21 op_table = {
22     "ADD": 0b0001, "AND": 0b0101, "NOT": 0b1001,
23     "LD": 0b0010, "LDI": 0b1010, "LDR": 0b0110,
24     "LEA": 0b1110, "ST": 0b0011, "STI": 0b1011, "STR": 0b0111,
25     "JMP": 0b1100, "JSR": 0b0100, "JSRR": 0b0100, "BR": 0b0000,
26     "TRAP": 0b1111,
27 }
28
29 trap_table = {
30     "GETC": 0x20, "OUT": 0x21, "PUTS": 0x22,
31     "IN": 0x23, "PUTSP": 0x24, "HALT": 0x25,
32 }
33
34 # ----- read input until .END -----
35 src = []
36 for line in sys.stdin:
37     line = line.rstrip("\n")
38     if line.strip().lower() == ".end":
39         break
40     src.append(line)
41
42 # ----- clean & normalize -----
43 clean = []
44 for line in src:
45     line = line.split(";")[0].strip()
46     if not line:

```

```
47     continue
48     line = re.sub(r"\s*,\s*", ", ", line)
49     clean.append(line)
50
51 # ----- first pass: labels -----
52 labels = {}
53 pc = None
54 lines = []
55
56 for line in clean:
57     parts_raw = line.split()
58     parts_up = [p.upper() for p in parts_raw]
59
60     # Handle .ORIG
61     if parts_up[0] == ".ORIG":
62         pc = to_int(parts_raw[1])
63         start = pc
64         continue
65
66     # Detect if the first token is a label
67     first = parts_raw[0]
68     op_candidate = parts_up[0]
69     if first.endswith(":") or (op_candidate not in op_table and op_candidate
70                               not in trap_table
71                               and not op_candidate.startswith(".") and
72                               op_candidate not in ["RET", "RTI"]):
73         label = first.rstrip(":")
74         labels[label] = pc
75         parts_raw = parts_raw[1:]
76         parts_up = parts_up[1:]
77
78     if not parts_raw:
79         continue
80
81     # Handle pseudo-ops to adjust PC
82     op = parts_up[0]
83     if op == ".BLKW":
84         n = to_int(parts_raw[1])
85         lines.append((pc, parts_raw, parts_up))
86         pc += n
87     elif op == ".STRINGZ":
88         s = " ".join(parts_raw[1:])[1:-1]
89         lines.append((pc, parts_raw, parts_up))
90         pc += len(s) + 1
91     else:
92         lines.append((pc, parts_raw, parts_up))
93         pc += 1
94
95 # ----- helper -----
```

```
94 def pcoff(pc_from, target, bits):
95     diff = target - (pc_from + 1)
96     max_val = (1 << bits)
97     if diff < 0:
98         diff = max_val + diff
99     return diff & (max_val - 1)
100
101 # ----- second pass -----
102 out = [to_bin16(start)]
103
104 def handle_instruction(op, parts_raw, parts_up, pc):
105     """Parse and handle different instruction types for LC-3 assembler."""
106     if op == ".FILL":
107         return [to_bin16(to_int(parts_raw[1]))]
108     if op == ".BLKW":
109         n = to_int(parts_raw[1])
110         return ["0000000000000000"] * n
111     if op == ".STRINGZ":
112         s = ".join(parts_raw[1:])[1:-1]"
113         s = s.replace("\n", "\n")
114         return [to_bin16(ord(ch)) for ch in s] + ["0000000000000000"]
115     if op in trap_table:
116         code = trap_table[op]
117         return [f"11110000{format(code, '08b')}"]
118     if op == "TRAP" and len(parts_raw) > 1:
119         imm = to_int(parts_raw[1]) & 0xFF
120         return [f"11110000{format(imm, '08b')}"]
121     if op == "RET":
122         return ["1100000111000000"]
123     if op == "RTI":
124         return ["1000000000000000"]
125     if op.startswith("BR"):
126         if len(op) == 2:
127             cond = "NZP"
128         else:
129             cond = op[2:]
130         n, z, p = ("N" in cond, "Z" in cond, "P" in cond)
131         label = parts_raw[1]
132         if label not in labels:
133             return []
134         target = labels[label]
135         imm9 = pcoff(pc, target, 9)
136         return [f"0000{int(n)}{int(z)}{int(p)}{format(imm9, '09b')}"]
137     if op == "JSR":
138         label = parts_raw[1]
139         if label not in labels:
140             return []
141         target = labels[label]
142         imm11 = pcoff(pc, target, 11)
```

```

143     return [f"01001{format(imm11, '011b')}"]
144 if op == "JSRR":
145     base = reg(parts_raw[1])
146     return [f"0100000{format(base, '03b')}000000"]
147 if op == "JMP":
148     base = reg(parts_raw[1])
149     return [f"1100000{format(base, '03b')}000000"]
150 if op in ["ADD", "AND"]:
151     rd = reg(parts_raw[1])
152     rs = reg(parts_raw[2])
153     immstr = parts_raw[3]
154     if immstr.startswith("#") or immstr.startswith("x"):
155         imm = to_int(immstr) & 0x1F
156         inst = (op_table[op] << 12) | (rd << 9) | (rs << 6) | 0x20 | imm
157     else:
158         rt = reg(immstr)
159         inst = (op_table[op] << 12) | (rd << 9) | (rs << 6) | rt
160     return [to_bin16(inst)]
161 if op == "NOT":
162     rd = reg(parts_raw[1])
163     rs = reg(parts_raw[2])
164     return [to_bin16((0b1001 << 12) | (rd << 9) | (rs << 6) | 0x3F)]
165 if op in ["LD", "ST", "LDI", "STI", "LEA"]:
166     opcode = op_table[op]
167     rd = reg(parts_raw[1])
168     label = parts_raw[2]
169     if label not in labels:
170         return []
171     target = labels[label]
172     imm9 = pcoff(pc, target, 9)
173     return [to_bin16((opcode << 12) | (rd << 9) | imm9)]
174 if op in ["LDR", "STR"]:
175     opcode = op_table[op]
176     dr = reg(parts_raw[1])
177     base = reg(parts_raw[2])
178     off = to_int(parts_raw[3]) & 0x3F
179     return [to_bin16((opcode << 12) | (dr << 9) | (base << 6) | off)]
180
181 for pc, parts_raw, parts_up in lines:
182     op = parts_up[0]
183     instructions = handle_instruction(op, parts_raw, parts_up, pc)
184     if instructions:
185         out.extend(instructions)
186
187 print("\n".join(out))

```

4.4.2 测试用例 example.rb

```
1 .ORIG x3000
2 ldi R1,GgPviBeMe
3 TRAP x7
4 JSRR R4
5 .fill x2
6 jmp R3
7 STR R5, R2, #1
8 .STRINGZ "wPTG12aiAAyPZ25BqAakH"
9 STI R1,3UEUAr
10 puts
11 ret
12 LDI r4, NdFRFwmZ
13 Ho4Yuix .STRINGZ "qlyxNYz0ex5sIHs4xEKOq"
14 OUT
15 LDI r3,TmB060D
16 HALT
17 jsrr R2
18 TRAP x-8
19 PUTS
20 IN
21 BR Ho4Yuix
22 HALT
23 JSRR r7
24 KVsyAq ld R1,ORKYrYrvy
25 TmB060D OUT
26 .BLKW 2
27 c_mA6EqFn and r5,R2,r5
28 st R3, 3UEUAr
29 trap x5
30 STI r6, KVsyAq
31 FU_p56KgRY
32 3UEUAr
33 NdFRFwmZ ST r1,ORKYrYrvy
34 ST R0, c_mA6EqFn
35 GgPviBeMe putsp
36 .blk 9
37 SviSJy1 out
38 RET
39 ldi r1,Ho4Yuix
40 and r5,R0,r7
41 JSRR R4
42 jsyHuX
43 McyiKg RTI
44 TRAP x-1
45 IN
46 JMP r0
47 BOOWxfWu
```

```
48 ORKYrYrvy .stringz "_QPM80juKHZhpNPUvjp3d"
49 RTI
50 ldr r4, r1, x-8
51 NOT r3,r3
52 add r5, r5, R2
53 LEA r4, aHwyACS
54 out
55 .FILL x-7
56 ret
57 LDR r3,R1,x-7
58 RTI
59 br c_mAK6EqFn
60 LD r4, KVsyAq
61 aHwyACS .FILL x2
62 EsOuVKrKC IN
63 ldr r6, R4, x1
64 ldr r1, R6, #5
65 jmp R6
66 rti
67 JSRR R5
68 trap #-6
69 sti r4, c_mAK6EqFn
70 br aHwyACS
71 ldi R0,jsyHuX
72 RTI
73 ST R7, FU_p56KgRY
74 STI R6, McyiKg
75 .END
```

4.4.3 测试用例 right.txt

```
1 0011000000000000
2 1010001001001001
3 111100000000111
4 0100000100000000
5 0000000000000010
6 1100000011000000
7 0111101010000001
8 000000001110111
9 0000000001010000
10 0000000001010100
11 0000000001000111
12 0000000001101100
13 000000000110010
14 0000000001100001
15 0000000001101001
16 0000000001000001
17 0000000001000001
18 0000000001111001
```

```
19 0000000001010000
20 0000000001011010
21 0000000000110010
22 0000000000110101
23 0000000001000010
24 0000000001110001
25 0000000001000001
26 0000000001100001
27 0000000001101011
28 0000000001001000
29 0000000000000000
30 1011001000101011
31 1111000000100010
32 1100000111000000
33 1010100000101000
34 0000000001110001
35 0000000001101100
36 0000000001111001
37 0000000001111000
38 0000000001001110
39 0000000001011001
40 0000000001111010
41 0000000001001111
42 0000000001100101
43 0000000001111000
44 0000000000110101
45 0000000001110011
46 0000000001001001
47 0000000001001000
48 0000000001110011
49 0000000000110100
50 0000000001111000
51 0000000001000101
52 0000000001011000
53 0000000001001111
54 0000000001110001
55 0000000000000000
56 1111000000100001
57 1010011000001001
58 1111000000100101
59 0100000100000000
60 11110000011111000
61 1111000000100010
62 1111000000100011
63 0000111111100010
64 1111000000100101
65 0100000111000000
66 0010001000011100
67 1111000000100001
```

```
68 0000000000000000  
69 0000000000000000  
70 0101101010000101  
71 0011011000000010  
72 1111000000000101  
73 1011110111111000  
74 0011001000010100  
75 0011000111111010  
76 1111000000100100  
77 0000000000000000  
78 0000000000000000  
79 0000000000000000  
80 0000000000000000  
81 0000000000000000  
82 0000000000000000  
83 0000000000000000  
84 0000000000000000  
85 0000000000000000  
86 1111000000100001  
87 1100000111000000  
88 1010001111001001  
89 0101101000000111  
90 0100000100000000  
91 1000000000000000  
92 1111000011111111  
93 1111000000100011  
94 1100000000000000  
95 0000000001011111  
96 0000000001010001  
97 0000000001010000  
98 0000000001001101  
99 0000000000111000  
100 0000000001001111  
101 0000000001101010  
102 0000000001110101  
103 0000000001001011  
104 0000000001001000  
105 0000000001011010  
106 0000000001101000  
107 0000000001110000  
108 0000000001001110  
109 0000000001010000  
110 0000000001010101  
111 0000000001110110  
112 0000000001101010  
113 0000000001110000  
114 0000000000110011  
115 0000000001100100  
116 0000000000000000
```

```
117 1000000000000000  
118 0110100001111000  
119 1001011011111111  
120 0001101101000010  
121 1110100000000111  
122 1111000000100001  
123 111111111111001  
124 1100000111000000  
125 0110011001111001  
126 1000000000000000  
127 000011111000110  
128 0010100111000001  
129 0000000000000010  
130 1111000000100011  
131 0110110100000001  
132 01100011100000101  
133 1100000110000000  
134 1000000000000000  
135 0100000101000000  
136 1111000011111010  
137 1011100110111100  
138 000011111110110  
139 1010000111001111  
140 1000000000000000  
141 001111110111100  
142 1011110111001100
```

5 实验测试与结果分析

5.1 测试用例

使用 `example.rb` 作为测试用例，该用例包含了 LC-3 的各类指令，如运算指令（ADD/AND/NOT）、访存指令（LDI/STR/STI）、控制转移指令（JMP/JSRR/BR）、陷阱指令（TRAP/GETC/OUT/PUTS）和伪指令（.ORIG/.FILL/.STRINGZ/.BLKW），同时包含多个标签和复杂的地址引用。（该大样例源自于程序提交后复制的测试点数据）

5.2 测试结果

运行 `main.py`，输入 `example.rb` 的汇编代码，输出的机器码与 `right.txt` 完全一致，验证了汇编器的正确性。以下给出了前 5 行的对比：

内存地址	汇编器输出（二进制）	正确结果（right.txt）	验证结果
x3000	0011000000000000	0011000000000000	一致
x3001	1010001001001001	1010001001001001	一致
x3002	1111000000000111	1111000000000111	一致
x3003	0110001000000000	0110001000000000	一致
x3004	1001001001001001	1001001001001001	一致

表 1: 测试结果对比表

5.3 结果分析

- 汇编器成功处理了各类指令的编码，包括寄存器操作数、立即数和标签地址的转换。
- 对于伪指令.STRINGZ，正确将字符串转换为 ASCII 码并在末尾添加 0；对于.BLKW，正确分配了指定数量的 0 填充内存空间。
- 基于 PC 的偏移量计算（如 BR、JSR 指令）和寄存器间接寻址（如 JSRR、LDR）处理正确，标签的地址解析无错误。
- 陷阱指令（如 OUT、PUTS、IN）和特殊指令（RET、RTI）的编码符合 LC-3 规范。

6 实验遇到的问题与解决方法

6.1 问题 1：负数偏移量的补码转换

问题描述：在计算 PC 偏移量时，负数偏移量的二进制表示不符合 LC-3 的补码要求，导致分支指令跳转错误。

解决方法：在 pcoff 函数中，通过计算偏移量的补码值 (`max_val + diff`, 其中 `max_val = 1 << bits`)，并与掩码 (`max_val - 1`) 按位与，确保偏移量为指定位数的无符号数。

6.2 问题 2：标签识别的滞后性

问题描述：一开始仅进行一遍扫描的时候，当标签定义在程序后面没有访问到的时候会视为不存在

解决方法：通过两遍扫描先将标签的存在性和位置确定再构建机器码

7 实验总结与体会

本次实验通过实现 LC-3 汇编器，深入理解了汇编器的工作原理和两趟扫描法的应用，掌握了 LC-3 指令的编码规则和地址计算方法。在实现过程中，遇到了数据类型转换、标签识别、偏移量计算等问题，通过分析 LC-3 体系结构规范和调试代码，逐一解决了这些问题，提升了编程能力和对计算机体系结构的理解。

同时，实验也发现了汇编器实现的可优化点，例如增加错误处理机制（如未定义标签、指令格式错误的提示）、支持更多伪指令和宏定义等。未来可以进一步完善汇编器的功能，使其更接近实际的汇编器。